



Databáze ◊ poznámky k přednášce

11. Fulltextové vyhledávání

verze z 9. prosince 2024

Pro motivaci si představme, že chceme zachytit část světa vymezenou námi vlastními filmy. Vlastníme tři filmy označené čísly jedna, dva a tři. U každého filmu nás zajímá pouze jeho název. První film se jmenuje *The Fellowship of the Ring*, druhý *The Two Towers* a třetí *The Return of the King*.

1 Indexy

Dokument je textový řetězec nad anglickou abecedou se znakem pro mezeru. Dokumenty budeme uzavírat mezi dvojité uvozovky. Například "The Fellowship of the Ring". **Token** je neprázdný textový řetězec nad malými písmeny anglické abecedy. Zapisujeme jej přímo, bez použití uvozek. Například `the` nebo `fellowship`.

Dokumentu D přiřadíme posloupnost **tokenů** t_1, \dots, t_n , kterou získáme tak, že znaky řetězce D převedeme na malá písmena a získaný řetězec rozdělíme na tokeny oddělené aspoň jednou mezerou. Například řetězec "The Fellowship of the Ring" rozdělíme na tokeny `the`, `fellowship`, `of`, `the`, `ring`. Zobrazení, které dokumentu přiřazuje posloupnost tokenů, označíme T . Tedy například:

$$T(\text{"The Fellowship of the Ring"}) = (\text{the}, \text{fellowship}, \text{of}, \text{the}, \text{ring})$$

Volněji budeme tokeny dokumentu ztotožňovat s jeho slovy.

Budeme používat systém Elasticsearch. Návod na jeho instalaci naleznete v souboru `11_elasticsearch.pdf`. Dotazem:

```
POST /_analyze
{
  "analyzer": "standard",
  "text": D
}
```

získáme tokeny dokumentu D :

```
{
  "tokens": [
    {
      "token":  $t_1$ ,
```

```

    },
    :
    {
      "token":  $t_n$ 
    }
  ]
}

```

Například:

```

POST /_analyze
{
  "analyzer": "standard",
  "text": "  The Fellowship    of the Ring"
}

{
  "tokens": [
    {
      "token": "the",
    },
    {
      "token": "fellowship",
    },
    {
      "token": "of",
    },
    {
      "token": "the",
    },
    {
      "token": "ring",
    }
  ]
}

```

Záznam je dvojice $\langle d, D \rangle$, kde d je přirozené číslo a D je dokument. Číslo d se nazývá **identifikátor** záznamu. Například $\langle 1, \text{"The Fellowship of the Ring"} \rangle$ je záznam s identifikátorem 1.

Konečná množina záznamů $\{\langle d_1, D_1 \rangle, \dots, \langle d_n, D_n \rangle\}$ s po dvou různými identifikátory se nazývá **index**. Například I_1 :

$$\{\langle 1, \text{"The Fellowship of the Ring"} \rangle, \langle 2, \text{"The Two Towers"} \rangle, \langle 3, \text{"The Return of the King"} \rangle\}$$

je index. **Vlastnost (záznamu)** je výroková forma $V(d, D)$ s dvěma proměnnými d a D , kde doména proměnné d je množina přirozených čísel a doména proměnné D je množina všech dokumentů. Například $V(d, D)$ může být: „Film d má název D .“ Říkáme, že index I má **charakteristickou vlastnost** $V(d, D)$, jestliže $\langle d, D \rangle \in I$, právě když $V(d, D)$. Také pak říkáme, že index I je **určen** vlastností $V(d, D)$. Například index I_1 má charakteristickou vlastnost: „Film d má název D .“ Platí, že $I = \{\langle d, D \rangle \mid V(d, D)\}$.

Prázdným indexem myslíme prázdnou množinu. Množinu všech identifikátorů $\{d_1, \dots, d_n\}$ záznamů v indexu I značíme $\text{Dom}(I)$. Například $\text{Dom}(I_1) = \{1, 2, 3\}$. Index I je zobrazení, které číslu z $\text{Dom}(I)$ přiřazuje dokument. Například $I_1(2) = \text{"The Two Towers"}$.

Indexová proměnná je proměnná, jejíž hodnota je index. Hodnotu indexové proměnné *index* označujeme I_{index} . Indexové proměnné *index* je přiřazena **charakteristická vlastnost** $V(d, D)$, což je vlastnost, která je vždy charakteristickou vlastností její hodnoty.

Indexovou proměnnou *index* deklaruujeme příkazem:

```
PUT /index
```

Hodnotou deklarované indexové proměnné je prázdný index. Například:

```
PUT /movie
```

Hodnota I_{movie} proměnné *movie* je prázdný index \emptyset . Charakteristická vlastnost indexové proměnné *movie* je: „Film d má název D .“

Pro index I a přirozené číslo d index

$$\{\langle d', D' \rangle \in I \mid d' \neq d\}$$

označíme $\text{Del}(I, d)$. Například pro výše uvedený index I_1 je $\text{Del}(I_1, 2)$ rovno:

$$\{\langle 1, \text{"The Fellowship of the Ring"} \rangle, \langle 3, \text{"The Return of the King"} \rangle\}$$

Pokud $d \notin \text{Dom}(I)$, pak $\text{Del}(I, d) = I$. Například $\text{Del}(I_1, 4) = I_1$.

Příkaz:

```
PUT /index/_doc/d
{
  "text": D
}
```

nastaví hodnotu proměnné *index* na $\text{Del}(I_{\text{index}}, d) \cup \{\langle d, D \rangle\}$.

Například:

```
PUT /movie/_doc/1
{
  "text": "The Fellowship of the Ring"
}
```

nastaví hodnotu proměnné `movie` na:

```
{⟨1, "The Fellowship of the Ring"⟩}
```

Po vykonání příkazů:

```
PUT /movie/_doc/2
{
  "text": "The Two Towers"
}
```

```
PUT /movie/_doc/3
{
  "text": "The Return of the King"
}
```

se hodnota proměnné `movie` změní na:

```
{⟨1, "The Fellowship of the Ring"⟩,
  ⟨2, "The Two Towers"⟩,
  ⟨3, "The Return of the King"⟩}
```

Příkaz:

```
DELETE /index/_doc/d
```

kde $d \in \text{Dom}(I_{index})$ nastaví hodnotu proměnné `index` na $\text{Del}(I_{index}, d)$. Například příkaz:

```
DELETE /movie/_doc/2
```

nastaví hodnotu proměnné `movie` na:

```
{⟨1, "The Fellowship of the Ring"⟩,
  ⟨3, "The Return of the King"⟩}
```

Záznam $\langle d, D \rangle$ zapíšeme řetězcem:

```
{
  "_id": d,
  "_source": {
    "text": D
  }
}
```

Například záznam $\langle 1, \text{"The Fellowship of the Ring"} \rangle$ zapíšeme řetězcem:

```
{
  "_id": 1,
  "_source": {
    "text": "The Fellowship of the Ring"
  }
}
```

Hodnota výrazu:

```
GET /index/_doc/d
```

kde $d \in \text{Dom}(I_{\text{index}})$ je $\langle d, I_{\text{index}}(d) \rangle$. Například:

```
GET /movie/_doc/1
```

```
{
  "_id": 1,
  "_source": {
    "text": "The Fellowship of the Ring"
  }
}
```

Index $\{\langle d_1, D_1 \rangle, \dots, \langle d_n, D_n \rangle\}$ zapíšeme řetězcem:

```
{
  "hits": {
    "hits": [
      <math>\langle d_1, D_1 \rangle</math>,
      :
      <math>\langle d_n, D_n \rangle</math>
    ]
  }
}
```

Například index:

```
{<1,"The Fellowship of the Ring">,  
<3,"The Return of the King">}
```

zapišeme řetězcem:

```
{  
  "hits": {  
    "hits": [  
      {  
        "_id": "1",  
        "_source": {  
          "text": "The Fellowship of the Ring"  
        }  
      },  
      {  
        "_id": "3",  
        "_source": {  
          "text": "The Return of the King"  
        }  
      }  
    ]  
  }  
}
```

Hodnota výrazu:

```
GET /index/_search  
{  
  "query": {  
    "match_all": { }  
  }  
}
```

je I_{index} , což je hodnota proměnné *index*. Například:

```
GET /movie/_search  
{  
  "query": {  
    "match_all": { }  
  }  
}  
  
{  
  "hits": {
```

```

    "hits": [
      {
        "_id": "1",
        "_source": {
          "text": "The Fellowship of the Ring"
        }
      },
      {
        "_id": "3",
        "_source": {
          "text": "The Return of the King"
        }
      }
    ]
  }
}

```

Příkaz:

```
DELETE /index
```

zruší proměnnou *index*. Například:

```
DELETE /movie
```

2 Podmínky

Zavedeme zobrazení Set , které posloupnosti termů (t_1, \dots, t_n) přiřadí množinu termů $\{t_1, \dots, t_n\}$, které se v posloupnosti vyskytují. Například:

$$\text{Set}(\text{the, fellowship, of, the, ring}) = \{\text{the, fellowship, of, ring}\}.$$

Podmínka je výraz, pro který lze u každého dokumentu rozhodnout, zda ji splňuje. Podmínce θ můžeme přiřadit její výrokovou formu $V_\theta(D)$ danou „Dokument D splňuje podmínku θ .“ Například podmínka **fellowship** má výrokovou formu „Dokument D obsahuje slovo **fellowship**.“.

Term t je podmínka, kterou dokument D splňuje, jestliže

$$t \in \text{Set}(T(D)).$$

Tedy dokument D splňuje t , jestliže se term t nachází mezi termy dokumentu D . Například dokument "The Fellowship of the Ring" splňuje podmínku **fellowship**.

Vezměme podmínku θ a index I . Pak index

$$\sigma_\theta(I) = \{ \langle d, D \rangle \in I \mid D \text{ splňuje } \theta \}$$

nazýváme **restrinkcí** I podle θ . Pokud $V_1(d, D)$ je charakteristická vlastnost indexu I , pak charakteristická vlastnost indexu $\sigma_\theta(I)$ je $V_2(d, D)$: „ $V_1(d, D)$ a $V_\theta(D)$.“
 Například pro index I_1 :

```
{⟨1, "The Fellowship of the Ring"⟩,
  ⟨2, "The Two Towers"⟩,
  ⟨3, "The Return of the King"⟩}
```

daný: „Film d má název D .“ je $\sigma_{\text{of}}(I_1)$ rovno

```
{⟨1, "The Fellowship of the Ring"⟩,
  ⟨3, "The Return of the King"⟩}.
```

a má charakteristickou vlastnost: „Film d má název D , který obsahuje slovo of.“

Vezměme podmínky $\theta_1, \dots, \theta_n$. Pak

$$\theta_1 \vee \dots \vee \theta_n$$

je podmínka, kterou dokument splňuje, jestliže splňuje aspoň jednu z podmínek $\theta_1, \dots, \theta_n$. Například $\sigma_{\text{two} \vee \text{king}}(I_1)$ je rovno:

```
{⟨2, "The Two Towers"⟩,
  ⟨3, "The Return of the King"⟩}.
```

a je dáno: „Film d má název D , který obsahuje slovo **two** nebo **king**.“

Vezměme dokument Q , pak řetězec:

```
{
  "match": {
    "text": Q
  }
}
```

zapisuje podmínku $q_1 \vee \dots \vee q_n$, kde $T(Q) = (q_1, \dots, q_n)$. Například řetězec

```
{
  "match": {
    "text": "Two King"
  }
}
```

zapisuje podmínku $\text{two} \vee \text{king}$.

Vezměme podmínku θ . Pak hodnotou výrazu:


```
GET /index/_search
{
  "query":  $\theta$ 
}
```

je restrikce $\sigma_{\theta}(I_{index})$. Například pro indexovou proměnnou `movie` danou: „Film d má název D .“ s hodnotou I_1 uvedenou výše je:

```
GET /movie/_search
{
  "query": {
    "match": {
      "text": "Two King"
    }
  }
}
```

rovno:

```
{
  "hits": {
    "hits": [
      {
        "_id": "2",
        "_source": {
          "text": "The Two Towers"
        }
      },
      {
        "_id": "3",
        "_source": {
          "text": "The Return of the King"
        }
      }
    ]
  }
}
```

s charakteristickou vlastností: „Film d má název D , který obsahuje slovo `two` nebo `king`.“

3 Kvalita splnění podmínky

Vezměme index I . Například index I_1 :

$$\{\langle 1, \text{"The Fellowship of the Ring"} \rangle, \\ \langle 2, \text{"The Two Towers"} \rangle, \\ \langle 3, \text{"The Return of the King"} \rangle\}.$$

Dále vezmeme term t . Pak číslo

$$n(I, t) = |\{\langle d, D \rangle \in I \mid t \in \text{Set}(T(D))\}|$$

udává **počet dokumentů** v I **obsahující term** t . Například $n(I_1, \text{the}) = 3, n(I_1, \text{of}) = 2, n(I_1, \text{king}) = 1, n(I_1, \text{three}) = 0$. **Počet záznamů v indexu** I označíme $N(I) = |I|$.

Číslo

$$IDF(I, t) = \ln \left(1 + \frac{N(I) - n(I, t) + 0.5}{n(I, t) + 0.5} \right)$$

udává **vzácnost termu** t v indexu I , kde IDF je zkratkou za *inverse document frequency*. Například:

$$IDF(I_1, \text{of}) = \ln \left(1 + \frac{3 - 2 + 0.5}{2 + 0.5} \right) \\ = \ln \left(\frac{4}{2.5} \right) = \ln(1.6) = 0,470\dots$$

Vezměme dokument D a označme $T(D) = (t_1, \dots, t_n)$. Pak číslo

$$f(t, D) = |\{i \mid t_i = t \text{ a } 1 \leq i \leq n\}|$$

je **počet výskytů termu** t v **dokumentu** D . Například:

$$f(\text{the}, \text{"The Return of the King"}) = 2.$$

Počet prvků posloupnosti termů $T(D) = (t_1, \dots, t_n)$ dokumentu D značíme $dl(D) = n$ a nazýváme **délka dokumentu** D . Například:

$$dl(\text{"The Return of the King"}) = 5.$$

Průměrná délka dokumentu v indexu I je dána vzorcem

$$avgdl(I) = \frac{\sum_i^n dl(D_i)}{N(I)}$$

kde $I = \{\langle d_1, D_1 \rangle, \dots, \langle d_n, D_n \rangle\}$. Například:

$$avgdl(I_1) = \frac{5 + 3 + 5}{3} = 4.\bar{3}$$

Frekvence termu t v dokumentu D vzhledem k indexu I je dána vzorcem

$$tf(I, t, D) = \frac{f(t, D)}{f(t, D) + k_1 \cdot \left(1 - b + b \cdot \frac{dl(D)}{avgdl(I)}\right)}$$

kde konstanta $k_1 = 1.2$ (*term saturation parameter*) a $b = 0.75$ (*length normalization parameter*). Například:

$$tf(I_1, \text{the, "The Return of the King"}) = \frac{2}{2 + 1.2 \cdot \left(1 - 0.75 + 0.75 \cdot \frac{5}{4.3}\right)} \approx 0.6$$

Definujeme zobrazení *score*, které přiřadí **kvalitu splnění** $score(I, D, \theta)$ podmínky θ dokumentem D z indexu I . Zobrazení definujeme zvlášť pro různé typy podmínek.

Pokud je podmínka term t , pak

$$score(I, D, t) = boost \cdot IDF(I, t) \cdot tf(I, t, D)$$

kde konstanta *boost* je $k_1 + 1 = 2.2$.

Pokud má podmínka tvar $\theta_1 \vee \dots \vee \theta_n$, pak

$$score(I, D, \theta_1 \vee \dots \vee \theta_n) = \sum_{i=1}^n score(I, D, \theta_i).$$

Definujeme zobrazení

$$\Sigma_\theta(I, I_G) = \{\langle\langle d, D \rangle, score(I_G, D, \theta)\rangle \mid \langle d, D \rangle \in I\},$$

kteří záznamu $\langle d, D \rangle \in I \subseteq I_G$ přiřadí kvalitu splnění podmínky θ dokumentem D z indexu I_G . Prvek zobrazení $\langle\langle d, D \rangle, s\rangle$, který se nazývá **zásah**, zapíšeme řetězcem:

```
{
  "_id": "d",
  "_score": s,
  "_source": {
    "text": D
  }
}
```

Celé zobrazení $\Sigma_\theta(I, I_G) = \{\langle\langle d_1, D_1 \rangle, s_1\rangle, \dots, \langle\langle d_n, D_n \rangle, s_n\rangle\}$ zapíšeme řetězcem:

```
{
  "hits": {
    "hits": [
      \langle\langle d_1, D_1 \rangle, s_1\rangle,

```

```

    :
    << $d_n, D_n$ >,  $s_n$ >
  ]
}
}

```

kde pro $1 \leq j, k \leq n$ platí, že pokud $s_j > s_k$, pak $j > k$. Zásahy jsou tedy uspořádány sestupně podle kvality.

Hodnotou výrazu:

```

GET /index/_search
{
  "query":  $\theta$ 
}

```

je zobrazení $\Sigma_{\theta}(\sigma_{\theta}(I_{index}), I_{index})$.

Například:

```

GET /movie/_search
{
  "query": {
    "match": {
      "text": "Two King"
    }
  }
}

{
  "hits": {
    "hits": [
      {
        "_id": "2",
        "_score": 1.1220688,
        "_source": {
          "text": "The Two Towers"
        }
      },
      {
        "_id": "3",
        "_score": 0.9227538,
        "_source": {
          "text": "The Return of the King"
        }
      }
    ]
  }
}

```

```
    ]
  }
}
```

Hodnotou výrazu:

```
GET /index/_search {
  "query":  $\theta$ ,
  "explain": true
}
```

navíc k zásahu přidáme položku `_explanation`, vyjadřující postup při výpočtu kvality zásahu, která má hodnotu:

```
{
  "value":  $score(I, D, \theta)$ ,
  "details": [
    {
      "details": [
        {
          "value":  $boost \cdot IDF(I, q_1) \cdot tf(I, q_1, D)$ ,
          "details": [
            { "value":  $boost$  },
            {
              "value":  $IDF(I, q_1)$ ,
              "details": [
                { "value":  $n(I, q_1)$  },
                { "value":  $N(I)$  }
              ]
            }, {
              "value":  $tf(I, q_1, D)$ ,
              "details": [
                { "value":  $f(t, D)$  },
                { "value":  $k_1$  },
                { "value":  $b$  },
                { "value":  $dl(D)$  },
                { "value":  $avgdl(I)$  }
              ]
            }
          ]
        }
      ]
    },
    :
  ]
}
```

Například:

```
GET /movie/_search
```

```
{
  "query": {
    "match": {
      "text": {
        "query": "Towers"
      }
    }
  },
  "explain": true
}

{
  "hits": {
    "hits": [
      {
        "_id": "2",
        "_score": 1.1220688,
        "_source": {
          "text": "The Two Towers"
        },
        "_explanation": {
          "value": 1.1220688,
          "details": [
            {
              "details": [
                {
                  "value": 1.1220688,
                  "details": [
                    { "value": 2.2 },
                    {
                      "value": 0.98082924,
                      "details": [
                        { "value": 1 },
                        { "value": 3 }
                      ]
                    }
                  ]
                },
                {
                  "value": 0.52000004,
                  "details": [
                    { "value": 1.0 },
                    { "value": 1.2 },
                    { "value": 0.75 },
                    { "value": 3.0 },
                    { "value": 4.3333335 }
                  ]
                }
              ]
            }
          ]
        }
      }
    ]
  }
}
```

```
    ]
  }
]
}
]
}
]
}
]
}
]
}
]
```

Otázky a úkoly na cvičení

1. Přiřaďte posloupnost tokenů k následujícím dokumentům.

- (a) The Life And Opinions Of Tristram Shandy
- (b) Emma
- (c) Nightmare Abbey
- (d) One Day in the Life of Ivan Denisovich
- (e) Life After Life

Porovnejte výsledky s Elasticsearch.

2. Vytvořte indexovou proměnnou `book` s charakteristickou vlastností „Kniha *d* má název *D*.“ a přidejte do ní dokumenty z předchozího úkolu.
3. Změňte název druhé knihy na `Frankenstein` a odstraňte z indexu knihu `Nightmare Abbey`.
4. Zjistěte hodnotu indexové proměnné `book`.
5. Určete hodnoty následujících restrikcí a ověřte výsledky systémem Elasticsearch.

- (a) $\sigma_{\text{life}}(I_{\text{book}})$
- (b) $\sigma_{\text{emma}}(I_{\text{book}})$
- (c) $\sigma_{\text{the\of}}(I_{\text{book}})$
- (d) $\sigma_{\text{fal\con\of\pinions\shandy}}(I_{\text{book}})$

Určete charakteristické vlastnosti získaných indexů.

6. Spočítejte následující hodnoty.

- (a) $n(I_{\text{book}}, \text{the})$

- (b) $n(I_{\text{book}}, \text{life})$
- (c) $N(I_{\text{book}})$
- (d) $IDF(I_{\text{book}}, \text{the})$
- (e) $IDF(I_{\text{book}}, \text{life})$
- (f) $f(\text{life}, \text{"Life After Life"})$
- (g) $f(\text{life}, \text{"Nightmare Abbey"})$
- (h) $dl(\text{"The Life And Opinions Of Tristram Shandy"})$
- (i) $avgdl(I_{\text{book}})$
- (j) $tf(I_{\text{book}}, \text{life}, \text{"Life After Life"})$
- (k) $tf(I_{\text{book}}, \text{life}, \text{"One Day in the Life of Ivan Denisovich"})$
- (l) $tf(I_{\text{book}}, \text{life}, \text{"The Life And Opinions Of Tristram Shandy"})$
- (m) $tf(I_{\text{book}}, \text{life}, \text{"Nightmare Abbey"})$
- (n) $score(I_{\text{book}}, \text{"Life After Life"}, \text{life})$
- (o) $score(I_{\text{book}}, \text{"One Day in the Life of Ivan Denisovich"}, \text{life})$
- (p) $score(I_{\text{book}}, \text{"One Day in the Life of Ivan Denisovich"}, \text{life} \vee \text{the})$

7. U následujících výrazů určete zásahy v jejich hodnotě a pořadí zásahů.

- (a) GET /book/_search


```
{
  "query": {
    "match": {
      "text": "Life"
    }
  }
}
```
- (b) GET /book/_search


```
{
  "query": {
    "match": {
      "text": "the"
    }
  }
}
```
- (c) GET /book/_search


```
{
  "query": {
    "match": {
      "text": "The LIFE"
    }
  }
}
```


Porovnejte vaše zásahy i pořadí s hodnotami výrazů vypočítanými systémem Elasticsearch. Nechte si systémem zobrazit výpočet kvality splnění zásahů.